

# SUPPORT DE COURS SI6



## EXPRESSIONS RÉGULIÈRES

date	révision
Mars 2018	Création
07/04/2018	Réorganisation des familles d'éléments (ancres, quantificateurs...) pour faciliter l'apprentissage
14/04/2018	Ajout d'exemples en JavaScript, Shell et PowerShell, PHP et Java
30/01/2019	Ajout d'une expression sur la date

## TABLE DES MATIÈRES

1	Introduction.....	3
2	Fonctionnement.....	3
3	La syntaxe REGEX.....	4
3.1	Les éléments de syntaxe.....	4
3.1.1	Les ancres (anchors).....	4
3.1.2	Les quantificateurs (quantifiers).....	4
3.1.3	Les classes de caractères.....	5
3.1.4	Les caractères (caractères de contrôle).....	5
3.1.5	Les groupes.....	5
3.2	Expressions simples.....	6
3.2.1	Recherche sur 're' et 'RE'.....	6
3.2.2	Recherche sur 'les' ou 'un'.....	6
3.2.3	Recherche sur les mots contenant 2 voyelles consécutives.....	6
3.2.4	Recherche sur les mots finissant par 're'.....	6
3.2.5	Recherche sur les mots commençant par 'la'.....	6
3.3	Expressions intermédiaires.....	7
3.3.1	Analyse d'une date.....	7
3.4	Expressions complexes.....	8
3.4.1	Valider la syntaxe d'une adresse de courrier.....	8
3.4.1.1	Expression 1.....	8
3.4.1.2	Expression 2.....	8
3.4.1.3	Expression 3.....	9
3.4.2	Valider la syntaxe d'un numéro de téléphone.....	10
3.4.3	Valider un mot de passe.....	10
4	Annexe.....	11
4.1	Regex Cheat Sheet.....	11
4.2	Regex en JavaScript.....	11
4.3	Regex en Shell linux.....	11
4.4	Regex en PowerShell (Microsoft).....	12
4.5	Regex en PHP.....	12
4.6	Regex en Java.....	12

## 1 INTRODUCTION

Les expressions régulières (aussi appelées REGEX pour REGular Exression) ne sont pas liées à un langage particulier mais au contraire, sont utilisées dans plusieurs langages.

En réalité, les expressions régulières décrivent un motif présent dans une chaîne de caractères et sont issues des langages formels développés en mathématiques dans les années 1940.

La lecture ou la rédaction d'une expression régulière dépend de sa complexité, mais son fonctionnement évite généralement le développement de nombreuses fonctions.

**Vous rencontrerez inévitablement des expressions régulières dans votre vie de développeur.**

## 2 FONCTIONNEMENT

Voici un exemple de REGEX (pris sur Wikipedia) qui permet de déterminer si une chaîne est une adresse de messagerie :

```
/[\w+.-]+@[ \w.-]+\.[a-zA-Z]{2,}/gi
```

La lecture serait une phrase du genre

Symbole REGEX	signification
[ ]	Définis un intervalle de caractères tel que...
\w	le motif de caractères alphanumériques de a à z, A à Z et 0 à 9
+	Pouvant se répéter (plusieurs occurrences)
.-	Et contenant aussi des tirets et des points
@	Puis un caractère arobase
etc	etc

Ce n'est clairement pas l'expression la plus simple mais cela donne rapidement une idée de la puissance de cette forme d'écriture !

Il faut pour cela utiliser des fonctions spéciales : en PHP, **preg\_filter()**, **preg\_grep()**, **preg\_match()**, etc.

Bien entendu, il existe des fonctions dans d'autres langages (Java ou JavaScript) et en C# :

```
var regex = new Regex(expression) ;
return regex.Matches(maChaîne) ;
```

Le symbole `\` a un rôle capital, car il permet d'échapper certains caractères (\*, ., +, { etc) mais rend également d'autres caractères spéciaux (\n, \r, \t...)

## 3 LA SYNTAXE REGEX

### 3.1 LES ÉLÉMENTS DE SYNTAXE

Les expressions régulières s'appuient sur un langage ayant plusieurs groupes d'éléments :

Ancre	Quantificateurs	Caractères	Groupes	Assertions	Modificateurs
Décrivent le début ou la fin d'un mot ou d'une chaîne.	Permettent la répétition un nombre de fois défini ou non.	Détermine les caractères spéciaux ayant un rôle particulier.	Délimiteurs facilitant la notion de groupe de symboles.	Capacité à répondre à une condition	Caractères hors de la recherche mais modifiant celle-ci
<i>Ex : ^signifie début de chaîne.</i>	<i>Ex : {2,5} indique une répétition de 2 à 5 fois max.</i>	<i>Ex : \s cible tous les espaces ensembles</i>	<i>Ex : (a b) indique un caractère comme a ou b.</i>	<i>Ex : ?! es une assertion négative</i>	<i>Ex : i est insensible à la casse.</i>

L'association de ces différents groupes permet des expressions complexes capables de valider si une chaîne de caractères contient un motif précis (adresse courriel, URL, date...)

#### 3.1.1 Les ancres (anchors)

Il s'agit de définir une position à partir de laquelle chercher le motif voulu. Les principales sont :

caractère	signification	exemple
^	Indique le début de la séquence de recherche (ligne)	^A accepte <b>A</b> rceau mais pas BANANA
\$	Indique la fin de la séquence de recherche (fin ligne)	A\$ accepte BANANA <b>A</b> mais pas Arceau
\b	Indique le bord d'un mot (à gauche = début, à droite = fin). <i>En effet, ^ et \$ s'applique à toute la chaîne.</i>	\b[D][\w]* accepte les mots comme <b>D</b> ans

#### 3.1.2 Les quantificateurs (quantifiers)

Ils permettent de définir un critère de répétition.

caractère	signification	exemple
*	Indique 0 ou plus	/fo*/ accepte <b>fo</b> otball et <b>fo</b> rt
+	Indique 1 ou plus	/ma+/ accepte <b>ma</b> ison et <b>ma</b> aison
?	0 ou 1 caractère seulement	/ba?c/ accepte <b>bc</b> <b>ba</b> c mais pas baac, barc
{ }	Selon le contenu des crochets : {2} pour 2 caractères uniquement, {2,} pour 2 et plus, {2,4} entre 2 et 4 caractères.	/w{3}/ accepte <b>www</b> mais pas ww ou <b>www</b> ww

### 3.1.3 Les classes de caractères

Généralement précédés du symbole `\`, ces caractères correspondent à un ensemble de symboles.

caractère	signification	exemple
<code>\w</code> <code>\W</code>	<code>\w</code> correspond à un caractère alphanumérique, <code>\W</code> ne doit pas correspondre à un caractère alphanu.	
<code>\d</code> <code>\D</code>	<code>\d</code> correspond à un symbole numérique (0 à 9) <code>\D</code> ne doit pas correspondre à un symbole numérique	<code>\d\d\d\d</code> accepte <b>2018</b> mais pas <code>1c21</code>
<code>\b</code> <code>\B</code>	<code>\b</code> correspond à un mot <code>\B</code> ne doit pas correspondre à un mot	
<code>\s</code> <code>\S</code>	<code>\s</code> correspond à tout type d'espace (plusieurs espaces consécutifs, tabulation...) <code>\S</code> ne doit pas correspondre à tout type d'espace.	
<code>.</code>	Correspond à n'importe quel symbole (une fois)	

### 3.1.4 Les caractères (caractères de contrôle)

L'association du symbole `\` avec ces symboles remplacent un caractère qui n'est pas accessible directement à l'utilisateur.

caractère	signification	exemple
<code>\n</code>	Un saut de ligne, caractère [Line Feed]	
<code>\r</code>	Un retour chariot, caractère [Carriage Return]	
<code>\t</code>	Une tabulation, caractère [ → ]	
<code>\x{N}</code>	Le caractère ayant pour code hexadécimal N	Le symbole € est <code>\x{20AC}</code>

### 3.1.5 Les groupes

Encadrés par les crochets [ et ], les symboles définissent un ensemble de valeurs autorisées ou non.

caractère	signification	exemple
<code>[0-7]</code>	Le symbole peut prendre les valeurs de 0 à 7	<code>[0-7]</code> Accepte <b>2</b> mais pas 8
<code>[abcdef]</code> <code>[^abc]</code>	Le symbole peut être une valeur entre a et f Le <code>^</code> est un complément : le symbole n'est pas a, b, c	
<code>[a-zA-Z]</code>	Le symbole est une lettre en minuscule ou majuscule	
<code>[0-9\-]</code>	Le symbole peut être un chiffre ou un signe moins	
<code>(a b)</code>	Le symbole peut être un a ou un b	<code>ch(at ien)</code> accepte <b>chat</b> ou <b>chien</b> . Pas <b>char</b>

Les groupes sont généralement suivis d'un quantificateur pour que le motif accepte plusieurs fois le groupe. Un numéro de téléphone à 10 chiffres sera `[0-9]{10}`.

## 3.2 EXPRESSIONS SIMPLES

Une expression régulière ne donne pas toute sa puissance si on recherche uniquement un mot.

Vous pouvez vous entraîner en ligne avec le site : <https://regex101.com/> ou <https://www.debuggex.com/>

Voici un morceau de texte que nous utiliserons pour nos tests :

Les expressions régulières (aussi appelées REGEX pour REGular Exression) ne sont pas liées à un langage particulier mais au contraire, sont utilisées dans plusieurs langages. En réalité, les expressions régulières décrivent un motif présent dans une chaîne de caractères et sont issues des langages formels développés en mathématiques dans les années 1940.

### 3.2.1 Recherche sur 're' et 'RE'

En programmation, il faudrait basculer le texte d'origine en majuscules ou minuscules puis faire une recherche sur 'RE' ou 're' selon le cas choisi.

L'expression régulière suivante règle le problème :

```
[rR][eE]
```

Cependant il existe aussi des modificateurs de motifs. L'expression suivante devrait fonctionner aussi :

```
/re/i
```

### 3.2.2 Recherche sur 'les' ou 'un'

On utilise le caractère | pour séparer les occurrences à accepter.

```
les|un
```

### 3.2.3 Recherche sur les mots contenant 2 voyelles consécutives

La solution est relativement simple : on commence par créer un intervalle ne contenant que des voyelles et on précise ensuite qu'il doit y avoir 2 occurrences consécutives.

```
[aeiou]{2}
```

### 3.2.4 Recherche sur les mots finissant par 're'

Le mauvais réflexe est d'utiliser le symbole \$ : il ne fonctionne que pour les fins de ligne. À la place :

```
re\b
```

### 3.2.5 Recherche sur les mots commençant par 'la'

Dans le même registre, le symbole ^ concerne le début de ligne. Pour les mots \b puis l'occurrence :

```
\bla
```

### 3.3 EXPRESSIONS INTERMÉDIAIRES

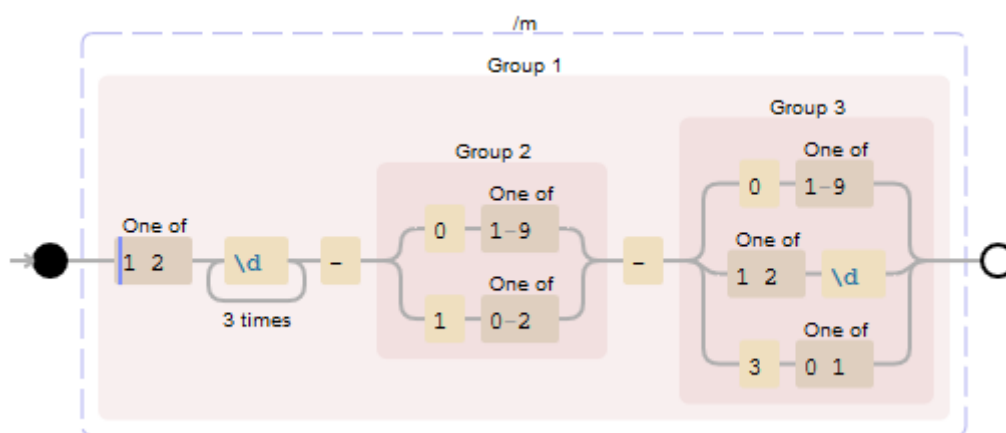
#### 3.3.1 Analyse d'une date

Voici le cas d'une date dans un format pratique pour les informaticiens :

YYYY-MM-DD

Si on veut tester une date de naissance ou actuelle, on part de l'année 1000 à 2999, du mois 01 à 12 et du jour 01 à 31.

Voici la présentation du raisonnement en REGEX :



et voici l'expression régulière :

```
([12]\d{3})-(0[1-9]|1[0-2])-(0[1-9]|12|\d3[01])
```

## 3.4 EXPRESSIONS COMPLEXES

Dans ce paragraphe, les expressions sont fournies "clé en main" : cela signifie qu'elles correspondent à des requêtes fréquentes. De manière classique en JavaScript, la vérification de champs de saisies.

### 3.4.1 Valider la syntaxe d'une adresse de courrier

Le format d'une adresse email est très reconnaissable au symbole @ entre deux chaînes de caractères. Malheureusement, aucune expression régulière ne donnera un résultat parfait.

#### 3.4.1.1 Expression 1

L'expression régulière de Wikipedia est intéressante mais simplifiée. Voici une autre expression<sup>1</sup>, compatible RFC 5322<sup>2</sup> un peu plus complète :

```
(?:[a-z0-9!#$%&'*/=\?^_`{|}~-]+(?:\.(?:[a-z0-9!#$%&'*/=\?^_`{|}~-]+)*)|"(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21\x23-\x5b\x5d-\x7f]|\[\x01-\x09\x0b\x0c\x0e-\x7f]*")@(?:[a-z0-9](?:[a-z0-9]*[a-z0-9])?\.)+[a-z0-9](?:[a-z0-9]*[a-z0-9])?|\[(?:(?25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9]?)\.){3}(?:25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9])?|[a-z0-9-]*[a-z0-9]:(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x53-\x7f]|\[\x01-\x09\x0b\x0c\x0e-\x7f])+\)])
```

REGULAR EXPRESSION

7 matches, 873 steps (~2ms)

```

/ (?: [a-z0-9!#$%&'*/=\?^_`{|}~-]+ (?: \. (?: [a-z0-9!#$%&'*/=\?^_`{|}~-]+ ) *) | "(?: [\x01-\x08\x0b\x0c\x0e-\x1f\x21\x23-\x5b\x5d-\x7f] | \[ \x01-\x09\x0b\x0c\x0e-\x7f ] * )" @ (?: [a-z0-9] (?: [a-z0-9]* [a-z0-9] ) ? \. ) + [a-z0-9] (?: [a-z0-9]* [a-z0-9] ) ? | \[ (?: (?: 25 [0-5] | 2 [0-4] [0-9] | [01]? [0-9] [0-9] ? ) \. ) {3} (?: 25 [0-5] | 2 [0-4] [0-9] | [01]? [0-9] [0-9] ? ) | [a-z0-9-]* [a-z0-9] : (?: [\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x53-\x7f] | \[ \x01-\x09\x0b\x0c\x0e-\x7f ] + ) \] ) ) /g

```

TEST STRING

SWITCH TO UNIT TESTS ▶

```

david@cf.fr
d@df.er.re.fr
d.avg.ret-exx01@fr05-ee.ffre.de
david.roumanet@ac-grenoble.comme
0dr@123fric.fr
david.roumanet+regex@ac-grenoble.fr
@ac-grenoble.fr
david@
david@ac-grenoble
d@e.
d@e.f

```

Le résultat reste imparfait, notamment le domaine racine (.com, .fr, etc.) mais exploitable.

#### 3.4.1.2 Expression 2

Si vous souhaitez une expression moins rébarbative :

```
\b([a-zA-Z0-9_-\.\+])@(\[[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3}\.|\([a-zA-Z0-9_-\.\+]\.+\))\b
```

1 <http://emailregex.com/>

2 RFC 5322 est la plus récente et rend obsolète les RFC 2822 et RFC 822.



Cette expression est relativement correcte. Notez également que l'expression doit être comprise entre les délimiteurs / et qu'il doit y avoir le modificateur /i pour ne pas être sensible à la casse.

Résultat :

REGULAR EXPRESSION v1 4 matches, 448 steps (~2ms)

```

/ \b([a-zA-Z0-9_-\.\+]+)@(\[\[0-9]{1,3}\. [0-9]{1,3}\. [0-9]{1,3}\. |(\[a-zA-Z0-9\-\.\+]\. )+)) /gi
([a-zA-Z]{2,4}|[0-9]{1,3})(\?)\b

```

TEST STRING SWITCH TO UNIT TESTS ▶

```

d@df.er.re.fr
d.avg.ret-exx01@fr05-ee.ffre.de
david.roumanet@ac-grenoble.comme
0dr@123fric.fr
david.roumanet+regex@ac-grenoble.fr
@ac-grenoble.fr
david@
david@ac-grenoble
d@e.
d@e.f

```

### 3.4.1.3 Expression 3

Dernière expression, la plus courte :

```
\b[A-Z0-9.\+_-]+@[A-Z0-9.-]+\.[A-Z0-9.-]+\b
```

et son résultat :

Regular Expression

```

/\b[A-Z0-9.\+_-]+@[A-Z0-9.-]+\.[A-Z0-9.-]+\b/ig

```

Test String

```

d@df.er.re.fr
d.avg.ret-exx01@fr05-ee.ffre.de
david.roumanet@ac-grenoble.comme
0dr@123fric.fr
david.roumanet+regex@ac-grenoble.fr
@ac-grenoble.fr
david@
david@ac-grenoble
d@e.
d@e.f

```

### 3.4.2 Valider la syntaxe d'un numéro de téléphone

Cette expression provient du site de Lucas WILLEMS<sup>3</sup> et permet de valider un téléphone en France.

```
(0|\+33)[1-9]( *[0-9]{2}){4}
```

Elle n'est cependant pas parfaite et nécessite de guider l'utilisateur dans la saisie.

Elle attend 0 ou +33 en début puis un chiffre. Ensuite elle accepte un espace et des chiffres par pair, 4 fois.

**REGULAR EXPRESSION**

```
/ (0|\+33)[1-9]( *[0-9]{2}){4}
```

**TEST STRING**

```
07.81.83.65.00
0781836500
07 81 83 65 00
+33 7 81 83 65 00
+33781836500
+337 81836500
```

### 3.4.3 Valider un mot de passe<sup>4</sup>

Cette expression<sup>5</sup> force un mot de passe ayant entre 8 et 15 caractères :

- au moins un caractère spécial \$ @ % \* + - \_ !
- au moins une minuscule
- au moins une majuscule
- au moins un chiffre

Regular Expression

```
/\b(?:[A-Z])(?:[a-z])(?:\d)(?:[-+!*$@%_])([-+!*$@%_ \w]{8,15})\b/i
```

Test String

```
azerty123
toto
LeschouxdeBruxelles$!
BC@racteres
Mcs1@ptesbuf
```

Voici l'expression codée :

```
\b(?:[A-Z])(?:[a-z])(?:\d)(?:[-+!*$@%_])([-+!*$@%_ \w]{8,15})\b
```

Cette expression introduit un concept intéressant de conditions, sous la forme

*(?=expression1)(?=expression2)* etc. :

prenez le temps de la décortiquer.

3 <https://www.lucaswillems.com/fr/articles/25/tutoriel-pour-maitriser-les-expressions-regulieres>

4 Suffisamment complexe

5 <http://codes-sources.commentcamarche.net/source/49715-validation-de-mot-de-passe>



## 4 ANNEXE

### 4.1 REGEX CHEAT SHEET

Anchors	Quantifiers	Character	Examples	POSIX	Groups	Modifiers
<b>Anchors</b> $\wedge$ Start of string, or start of line in multi-line pattern $\backslash A$ Start of string $\$$ End of string, or end of line in multi-line pattern $\backslash Z$ End of string $\backslash b$ Word boundary $\backslash B$ Not word boundary $\backslash k$ Start of word $\backslash b$ End of word	<b>Quantifiers</b> $*$ 0 or more $+$ 1 or more $?$ 0 or 1 $\{3\}$ Exactly 3 $\{3,5\}$ 3 or more $\{3,4,5\}$ 3, 4 or 5  <b>String</b> <b>Replacement</b> $\$n$ nth non-capturing group $\$2$ "xyz" in $P(abc-(xy-z))S/$ $\$1$ "xyz" in $P(?(a-bc)(xyz)S/$ $\$$ Before matched string $\$'$ After matched string $\$+$ Last matched string $\$&$ Entire matched string	<b>Character Classes</b> $\backslash c$ Control character $\backslash s$ White space $\backslash S$ Not white space $\backslash d$ Digit $\backslash D$ Not digit $\backslash w$ Word $\backslash W$ Not word $\backslash x$ Hexade-cimal digit $\backslash O$ Octal digit  <b>Special</b> $\backslash n$ New line $\backslash r$ Carriage return $\backslash t$ Tab $\backslash v$ Vertical tab $\backslash f$ Form feed $\backslash oox$ Escal character xxx $\backslash xhh$ Hex character hh	<b>Metacharacter</b> $\wedge abc$ $abc\$$ $abc.endsinabc.123abc...$ $a.c$ $abc.aac.aac.adc.aec...$ $bill ted$ $ted bill$ $ab{2}c$ $abc$ $a b c$ $abc.aBc$ $(abc){2}$ $abcabc$ $ab*c$ $abc.abc.abc.abbcc...$ $ab+c$ $abc.abc.abbcc...$ $ab?c$ $abc.abc$ $a c$ <b>Sample</b> $([A-Za-z0-9-]+)$ Letters, numbers and hyphens $(\d{1,2})\d{1,2}\d{1,2}$ Phone no. 12345678 $(([s]+)?= (p q f p ng))\d{2}$ jpg, gif or png image $(^[1-9]{1})\$(^[1-4]{1}){0-9}{1}$ $\$(^50\$)$ Any number from 1 to 50 inclus $(#?[A-Fa-f0-9-]{3}){((A-Fa-f0-9-){3})?}$ Hexadecimal colour code $((?=[^d](?=[a-z])?=[A-Z]) (8,15))$ 8 to 15 character string with at least one upper case letter, one lower case letter, and one digit (useful for passwords) $(w+@[a-zA-Z_]+?\.[a-zA-Z]{2,6})$ Email addresses $(\<([^\>]+)\>)$ HTML Tag	<b>POSIX</b> $[upper]$ Upper case letters $[lower]$ Lower case letters $[alpha]$ All letters $[alnum]$ Digits and letters $[digit]$ Digits $[xdigit]$ Hexade-cimal digits $[punct]$ Punctuation $[cntrl]$ Control characters $[print]$ Printed characters $[graph]$ Printed characters and spaces $[word]$ Digits, letters and underscore	<b>Groups and Ranges</b> $.$ Any character except new line $(n)$ Group $(a b)$ a or b $(...)$ Group $(?...)$ Passive (non-capturing) group $[abc]$ Range (a or b or c) $[^abc]$ Not a or b or c $[a-q]$ Letter from a to q $[A-Q]$ Upper case letter from A to Q $[0-7]$ Digits from 0 to 7 $\backslash n$ nth group/sub-pattern	<b>Modifiers</b> $g$ Global match $i$ Case-insensitive $m$ Multiple lines $s$ Treat string as single line $x$ Allow comments and white space in pattern $e$ Evaluate replacement $U$ Ungreedy pattern  <b>Assertions</b> $?$ Lookahead assertion $?!$ Negative lookahead $?<$ Lookbehind assertion $?=$ or $?<$ Negative lookbehind $?>$ Once-only Subexp-ression $?0$ Condition [if then] $?0 $ Condition [if then else] $?#$ Comment

### 4.2 REGEX EN JAVASCRIPT

JavaScript propose une classe RegExp qui fonctionne comme suit :

```
function checkZorg(chaine) {
    var motif = /[M-Z]or.$/g ;
    return motif.test(chaine) ;
}
```

### 4.3 REGEX EN SHELL LINUX

L'usage de commande ou de script en Shell permet d'utiliser les expressions régulières. Elles sont alors précédées du symbole ~ :

```
nom="Zorg"
if [[ $nom =~ [M-Z]or.$ ]] ; then
    echo vrai
fi
```

## 4.4 REGEX EN POWERSHELL (MICROSOFT)

La commande -match utilise les expressions régulières :

```
$nom = "Zorg"
if ($nom -match "[M-Z]or.$") {
    write-host "il y a correspondance"
}
```

A noter : en PowerShell, le caractère d'échappement est ` : il s'obtient avec [Alt Gr] + [7]

## 4.5 REGEX EN PHP

La commande -match utilise les expressions régulières :

```
<?php
if (preg_match("/[M-Z]or.$/g", "Zorg")) {
    echo "A match was found.";
} else {
    echo "A match was not found.";
}
?>
```

## 4.6 REGEX EN JAVA

Voici un code complet :

```
import java.io.*;
import java.util.regex.*;

public class testRegex {
    private static Pattern motif;
    private static Matcher matcher;

    public static void main(String args[]) {
        motif = Pattern.compile("[M-Z]or.$");
        matcher = pattern.matcher("Zorg mon ami Zorg");
        while(matcher.find()) {
            System.out.println("Trouvé !");
        }
    }
}
```